

This application is submitted in the name of the following inventors:

<u>Inventor</u>	<u>Citizenship</u>	<u>Residence City and State</u>
Blake LEWIS	United States	Palo Alto, California
Kayuri PATEL	Kenya	Cupertino, California
Ray CHEN	United States	Campbell, California

The assignee is Network Appliance, Inc., a California corporation having an office at 495 East Java Dr., Sunnyvale, CA 94089.

Title of the Invention

Reserving File System Blocks

Background of the Invention

1. Field of the Invention

This invention relates to management of a file system for a file server. In particular, the invention concerns reserving unallocated blocks of the file system based upon a file size for a file.

2. Description of the Related Art

Some conventional file servers (also called "filers") manage space in a file system by allocating blocks to a file as data for that file is written to the file system. Thus, if the file system runs out of space, a file could end up partially written.

1 Other conventional file servers, in particular those running CIFS, can allocate
2 blocks for an entire file upon creation of the file. These file servers write zero data to the file
3 system for these blocks. However, this approach is expensive in terms of time required to write
4 the zero data. Furthermore, this approach is actually counterproductive for file servers that use a
5 write anywhere file system layout, also known as WAFL file servers.

6
7 In WAFL file servers, when a file is overwritten, new data is written to new
8 blocks, and then the previously allocated blocks are released. Thus, new data is written to
9 different blocks than the previously allocated blocks, resulting in use of extra space for the new
10 blocks until the previously allocated blocks are released. If the file server is close to full, this
11 duplication of blocks could use up the remaining blocks, preventing complete writing of the data.

12 //

13 //

14 //

15 //

Summary of the Invention

Accordingly, what is needed is a way to ensure that enough blocks are reserved for a file so as to ensure that the entire file can be written to a file system, without actually allocating disk blocks to the file.

In one aspect, the invention addresses the foregoing need through a method of managing a file system for a file server. According to the method, a file operation is received that signals a reservation operation for a file having a file size. Preferably, the file system uses a write anywhere file system layout, the file operation that signals the reservation operation is a zero length write request, and the file operation that signals the reservation operation includes a parameter that specifies the file size. A number of blocks needed to be reserved to accommodate the file is computed. Preferably, computing the number of blocks needed to be reserved to accommodate the file includes determining a total number of direct and indirect blocks needed to accommodate the file size, and subtracting a total number of blocks already allocated for the file and a total number of cached unallocated blocks for the file from the total number of direct and indirect blocks needed to accommodate the file size. Unallocated blocks are reserved in the file system, with the number of reserved blocks equal to the number of blocks needed to be reserved to accommodate the file. Reserving the number of blocks preferably includes setting a flag in an inode for the file that indicates blocks have been reserved for the file, and incrementing a reserved block count in a file system information block by the number of blocks needed. The reserved block count indicates how many unallocated blocks have been reserved for files in the file system.

In a preferred embodiment, the method also includes the step of checking that a number of available blocks in the file system is greater than the number of blocks needed to be

1 reserved to accommodate the file. An error is returned in a case that the number of available
2 blocks is less than the number of blocks needed. The number of available blocks in the file
3 system preferably is determined by subtracting a number of allocated blocks, a number of cached
4 unallocated blocks (i.e., delayed allocated blocks), and a number of reserved blocks from a total
5 number of blocks in the file system, and adding to this a number of reserved cached unallocated
6 blocks.

7
8 Also in the preferred embodiment, the file server checks that the number of blocks
9 needed to be reserved to accommodate the file does not exceed a remainder of a quota for an
10 owner of the file. An error is returned in a case that the number of blocks needed exceeds the
11 remainder of the quota.

12
13 When data is written to the file system, blocks are cached in a buffer cache. At a
14 later point in time, the blocks are stored to storage for the file system. Prior to writing the blocks
15 to disk, the blocks are actually allocated to the files. Reservation of those blocks is released as
16 the blocks are written to storage. Releasing reservation of blocks is accomplished by
17 decrementing the reserved block count in the file system information block by a number of
18 released blocks.

19
20 By virtue of the foregoing operations, a file server reserves unallocated blocks for
21 a file for which file reservation semantics are activated. These reserved blocks are not actually
22 allocated by the reservation process. Rather, a count is maintained of how many blocks need to
23 be kept available for the file system. This count is utilized when space availability for the file
24 system is checked, thereby helping to ensure that enough blocks are available for files that have
25 reserved file space.

1 In another aspect, the foregoing method handles receipt of a file operation that
2 signals a reservation operation for a file for which reservation has already been performed, in
3 which the reservation operation specifies a new file size different from a current file size for the
4 file. When such a file operation is received, additional blocks may need to be reserved to
5 accommodate the new file. According to the method, a current file size for the file and the new
6 file size are compared. In a case that the current file size exceeds the new file size, remaining
7 block reservations for the file are released. In a case that the new file size exceeds the current file
8 size, an additional number of unallocated blocks are reserved in the file system. According to
9 this embodiment of the invention, the additional number of unallocated blocks to be reserved
10 equals a difference between a total number of direct and indirect blocks required by the new file
11 size and a total number of direct and indirect blocks required by the current file size.

12 By virtue of the foregoing operation, changes in file size for a file can be
13 appropriately reflected in the block reservations.

14 Each of the foregoing methods can be used in conjunction with the others in
15 various combinations to perform reservation operations. The invention also can be embodied in
16 apparatuses such as file servers and/or other hardware configured to perform the foregoing
17 methods, computer readable code by itself or embodied in a computer program product to cause a
18 computer to perform the foregoing methods, and a memory storing information including
19 instructions executable by a processor to perform the foregoing methods.
20

21
22 This brief summary has been provided so that the nature of the invention may be
23 understood quickly. A more complete understanding of the invention may be obtained by
24 reference to the following description of the preferred embodiments thereof in connection with
25 the attached drawings.
26

Brief Description of the Drawings

Figure 1 shows an environment in which a file server manages a file system according to the invention.

Figure 2 is a block diagram of an organization of a file system in which blocks can be reserved for files.

Figure 3 is a flowchart for explaining reservation according to the invention of blocks for a file.

Figure 4 is a flowchart for explaining computation of a number of blocks needed for a file.

Figure 5 is a flowchart for explaining determination of how many blocks are available for a file system.

Figure 6 is a flowchart for explaining reservation of unallocated blocks for a file.

Figure 7 is a flowchart for explaining allocation and writing of data to a file in a file system in which blocks can be reserved for files.

Figure 8 is a flowchart for explaining block reservation when a file size changes.

1 Description of the Preferred Embodiment

2

3 In the following description, a preferred embodiment of the invention is described

4 with regard to preferred process steps and data structures. However, those skilled in the art

5 would recognize, after perusal of this application, that embodiments of the invention may be

6 implemented using one or more general purpose processors or special purpose processors

7 adapted to particular process steps and data structures operating under program control, that such

8 process steps and data structures can be embodied as information stored in or transmitted to and

9 from memories (e.g., fixed memories such as DRAMs, SRAMs, hard disks, caches, etc., and

10 removable memories such as floppy disks, CD-ROMs, data tapes, etc.) including instructions

11 executable by such processors (e.g., object code that is directly executable, source code that is

12 executable after compilation, code that is executable through interpretation, etc.), and that

13 implementation of the preferred process steps and data structures described herein using such

14 equipment would not require undue experimentation or further invention.

15

16 Fig. 1 shows an environment in which a file server manages a file system

17 according to the invention. In Fig. 1, clients 1 send data to and retrieve data from server 2, which

18 stores the data in file system 3. File system 3 includes storage 4 such as multiple hard disk

19 drives, multiple optical drives, or any other mass storage. Information is stored in storage 4 in

20 blocks, as explained in more detail below with reference to Fig. 2. Preferably, file system 3 uses

21 a write anywhere file system layout (WAFL).

22 When data is sent to server 2, that data is initially stored in buffer cache 5 before

23 being written to file system 3. Data in buffer cache 5 also is stored in blocks, which are shown as

24 cached blocks 6 in Fig. 1. According to the invention, each cached block 6 preferably includes

25 flag 7 that indicates whether or not the block was cached in buffer cache 5 after block reservation

26 according to the invention was activated for the file to which the block belongs. The operation of

1 flag 7 is discussed below with respect to Figs. 3 to 8.

2
3 Blocks in buffer cache 5 can be copies of allocated blocks already written to
4 storage 4, or unallocated blocks which have not yet been written to storage 4. An allocated block
5 has a block identification number (not shown) associated therewith that indicates to which block
6 of storage 4 it corresponds. An unallocated blocks has no such block identification number
7 associated therewith, or has a block identification number of zero.

8
9 When a new block is created in buffer cache 5, that block is "clean." When data
10 is written to the block, the block is "dirtied." File server 2 preferably uses delayed allocation for
11 these dirty blocks. According to delayed allocation, blocks in storage 4 are not immediately
12 allocated for the dirty blocks, and the dirty blocks preferably are not immediately written to
13 storage 4. Thus, the dirty blocks are cached unallocated blocks. Such cached unallocated blocks
14 are also called delayed allocated blocks.

15
16 Actual writing of the cached unallocated blocks (i.e., delayed allocated blocks)
17 occurs at certain predetermined times or after the occurrence of certain predetermined conditions
18 (i.e., NVRAM for persistent storage of file server operations fills up, a preset number of buffers
19 or blocks are cached, etc.). Then, blocks in storage 4 are allocated for dirtied blocks in buffer
20 cache 5, and the dirtied blocks are actually stored in storage 4.

21
22 File server 2 also includes incore information structure 8 with incore inodes 9 and
23 incore information 10. An incore inode 9 is present for each file for which blocks are stored in
24 buffer cache 5. The incore inodes associate the blocks with the file in a manner similar to inode
25 14 discussed below with respect to Fig. 2.

1 Each inode also includes a count of the number of cached unallocated
2 blocks (i.e., delayed allocated blocks) associated with the inode. Incore information 10
3 preferably includes a count of how many total cached unallocated blocks are stored in buffer
4 cache 5.

5
6 Fig. 2 is a block diagram of an organization of file system 3. For the sake of
7 brevity, many of the conventional features of file system 3 are not illustrated in Fig. 2 and are not
8 discussed herein. Those skilled in the art will appreciate that a great many ways exist to
9 implement these conventional features without departing from the invention.

10
11 File system 3 includes at least one system information block 11. This system
12 information block includes at least block allocation information 12 and reserved block count 13
13 according to the invention. Block allocation information 12 includes at least information from
14 which can be determined how many blocks of file system 3 are allocated to files. Reserved block
15 count 13 includes a count of how many blocks in file system 3 have been reserved according to
16 the invention.

17
18 An inode such as inode 14 is associated with each file stored in file system 3.
19 Inode 14 is stored in storage 4. Inode 14 includes at least flag 15 that indicates whether or not
20 block reservation according to the invention is active for the file to which the inode belongs. The
21 operation of flag 15 is discussed below with respect to Figs. 3 to 8. Inode 14 also can include file
22 size information for its associated file.

23
24 Inode 14 for a file further includes information associating that file with blocks 16
25 in file system 3. These blocks preferably are 4 kilobytes long. Blocks 16 can be direct blocks, in
26 which case data for the file is stored directly in the blocks. Blocks 16 also can be indirect blocks

1 that store information associating other blocks with the file. Those other blocks in turn can store
2 data or can point to yet more blocks.

3
4 For a file size of less than 64 bytes, no blocks are needed. Instead, data for the file
5 is stored directly in inode 14. In a typical block configuration, each inode can point to up to
6 sixteen other blocks. Thus, for file sizes between 64 bytes and 64 kilobytes, up to sixteen direct
7 blocks are utilized. Each block in turn can point to 1024 other blocks. Thus, one level of
8 indirection can accommodate files between 64 kilobytes and 64 megabytes (16 blocks pointing to
9 1024 blocks of 4 kilobytes each). Two levels of indirection can accommodate file sizes up to 64
10 gigabytes (16 blocks pointing to 1024 blocks, which each in turn point to 1024 blocks of 4
11 kilobytes each). More levels of indirection can be utilized, as needed. Of course, the invention is
12 equally applicable to file systems that utilize different size blocks than 4 kilobytes and that are
13 organized in different arrangements than shown in Fig. 2

14
15 As data is written to a file in file system 3, the data is cached in block 6 in the
16 buffer cache. Prior to writing the block 6 to storage 4, blocks are allocated to that file. File server
17 2 manages the storage of the data in the blocks and keeps track of block allocation. According to
18 the invention, file server 2 also can reserve blocks for files before data is written to those files.
19 File server 2 uses reserved block count 13 to keep track of how many blocks are reserved.

20
21 Reserved blocks preferably are not actually assigned to a file. Instead, file server
22 2 ensures that at least the number of blocks indicated by reserved block count 13 are kept free.
23 Then as data is written to a file for which reservation is activated, the data is cached in block 6 in
24 the buffer cache. Prior to writing the block 6 to storage 4, blocks are allocated to that file.
25 Because those blocks are now allocated, space for the blocks no longer needs to be reserved in
26 file system 3, so reserved block count 13 is decremented accordingly. These operations are

discussed in more detail below.

Block allocation information 12, reserved block count 13, and inodes 14 with flags 15 preferably are stored in file system 3 on storage 4. Thus, whenever a snapshot of file system 3 is made for backup purposes, sufficient block allocation and reservation information is stored in the snapshot to reconstruct the block allocation and reservation status from the snapshot.

Fig. 3 is a flowchart for explaining reservation according to the invention of blocks for a file.

Briefly, the invention manages a file system for a file server. A file operation is received that signals a reservation operation for a file having a file size. A number of blocks that need to be reserved to accommodate the file is computed. A number of unallocated blocks is reserved in the file system equal to the number of blocks needed to be reserved to accommodate the file.

In more detail, in step S301, file server 2 receives a file operation that signals a reservation operation for a file. In the preferred embodiment of the invention, this file operation is a zero length write request. This write request preferably includes a parameter that identifies a size for the file for which the reservation is to be made. Alternatively, the file size can already be set in file system 3, possibly through a previous write command.

In step S302, a number of blocks needed to be reserved to accommodate the file is computed. This step is shown in more detail in Fig. 4.

1 In step S401 of Fig. 4, a number of blocks needed to accommodate the file size is
2 computed. This number preferably is irrespective of information such as actual block allocation
3 and reservation data. The number of blocks includes both direct blocks and indirect blocks, as
4 described above with respect to Fig. 2.

5
6 For example, for a file of size 64 kilobytes, a preferred embodiment of the
7 invention needs 16 direct blocks and no indirect blocks, for a total of 16 blocks. The number
8 computed in step S401 preferably is generated by an algorithm that takes a file size as an input
9 and generates a number of blocks as an output.

10
11 In step S402 of Fig. 4, a number of blocks already allocated to the file, for
12 example as a result of previous write operations, is subtracted from the number of blocks needed
13 to accommodate the file size. Information about the number of blocks allocated to the file
14 preferably is retrieved or derived from block allocation information 12 and/or inodes 14.

15
16 In step S403, a number of cached unallocated blocks for the file is subtracted from
17 the result of step S402. While these blocks are unallocated, space for them is already accounted
18 for through delayed allocation counters stored in the file server's incore structure, so there is no
19 need to include these blocks in the reservation count. The number of unallocated cached blocks
20 (i.e., delayed allocated blocks) for the file is stored in the associated incore inode 9, and the total
21 count of delayed allocated blocks in buffer cache 5 is stored in incore information 10.

22
23 Thus, steps S401 through S403 determine a total number of direct and indirect
24 blocks needed to accommodate the file size, and subtract from this number a total number of
25 blocks already allocated for the file and a total number of cached unallocated blocks for the file.
26 The result is the number of blocks that need to be reserved to accommodate the file. Of course,

1 steps S402 and S403 can occur simultaneously or in a different order, as long as the number of
2 needed blocks is determined.

3
4 Returning to Fig. 3, the number of blocks available for file system 3 is determined
5 in step S303. This step is shown in more detail in Fig. 5.

6
7 In step S501 of Fig. 5, a total number of blocks for file system 3 is determined.
8 This number is a substantially fixed number that usually only changes when storage 4 is
9 physically altered, for example by permanent corruption of blocks (i.e., bad blocks), changes in
10 hardware (i.e., addition or removal of disk drives), reformatting, and the like.

11 In step S502, a number of allocated blocks is subtracted from the total number of
12 blocks. The number of allocated blocks preferably is determined from block allocation
13 information 12 in system information block 11 of file system 3.

14 In step S503, a number of cached unallocated blocks (i.e., delayed allocated
15 blocks) is subtracted from the result of step S503. The number of delayed allocated blocks is
16 determined from incore information 10 for file system 3.

17
18 In step S504, a number of reserved blocks is subtracted from the result of step
19 S503. The number of reserved blocks preferably is determined from reserved block count 13 in
20 system information block 11 of file system 3.

21
22 In step S505, a number of reserved cached unallocated blocks is added to the
23 result of step S504. File server 2 preferably keeps track of how many cached unallocated blocks
24 in buffer cache 5 have flag 7 set, thereby indicating that those blocks are blocks for which
25 reservation according to the invention is activated. This information preferably is kept up to date
26

1 in incore information 10 for file system 3. This step is necessary since the reserved cached
2 unallocated blocks are accounted for in the delayed allocation counters in incore information 10,
3 since actual blocks in storage 4 have not yet been allocated for these blocks.

4
5 The result of steps S501 through S505 is a number of available blocks for the file
6 system. Of course, steps S501 through S505 can occur simultaneously or in a different order, as
7 long as the number of available blocks is determined.

8
9 Returning to Fig. 3, in step S304, the number of blocks needed to be reserved is
10 compare with the number of available blocks. If the number of needed blocks exceeds the
11 number of available blocks, not enough space is available in file system 3, and flow is diverted to
12 step S305, where an error is returned. Otherwise, flow proceeds to step S306.

13
14 In step S306, a number of blocks remaining in the file owner's quota is
15 determined. This operation is only performed if quotas are being enforced for file system 3.
16 Then, in step S307, the quota is compared against the number of blocks needed to be reserved. If
17 not enough blocks remain in the quota, flow is diverted to step S308, where an error is returned.
18 Otherwise, flow proceeds to step S309.

19
20 Unallocated blocks are reserved in step S309. The number of blocks reserved is
21 equal to the number of blocks needed, as determined in step S302. Step S309 is shown in more
22 detail in Fig. 6.

23
24 In step S601, flag 15 in inode 14 for the file is set. This flag indicates that
25 reservation semantics are being enforced for the file. Then, in step S602, reserved block count
26 13 in system information block 11 is incremented by the number of blocks needed to be reserved

1 for the file.

2
3 After step S309, blocks have been reserved for the file.

4
5 Fig. 7 is a flowchart for explaining allocation and writing of data to a file in a file
6 system in which blocks can be reserved for files.

7
8 When data is to be written to file system 3 for a file, step S701 checks flag 15 of
9 inode 14 for that file to determine whether or not block reservation has been performed for the
10 file. If reservation has been performed, flow proceeds to step S705. Otherwise, flow proceeds to
11 step S702 for space availability checking.

12
13 The number of blocks available for file system 3 is determined in step S702.
14 Preferably, step S702 subtracts a number of allocated blocks, a number of cached unallocated
15 blocks (i.e., delayed allocated blocks), and a number of reserved blocks from a total number of
16 blocks in the file system, and adds to this a number of reserved cached unallocated blocks.

17
18 In step S703, the number of blocks needed to be allocated is compared with the
19 number of available blocks. If the number of needed blocks exceeds the number of available
20 blocks, not enough space is available in file system 3, and flow is diverted to step S704 where an
21 error is returned. Otherwise, flow proceeds to steps S705.

22
23 In step S705, the data is written to a block or blocks in buffer cache 5 for file
24 system 3. If reservation has been activated for the file, file server 3 ensures that flags 7 are set for
25 the blocks.
26

1 Flow proceeds from step S705 to step S706 in order to write data from buffer
2 cache 5 to storage 4. As mentioned above, data preferably is written from the buffer cache to the
3 storage after the occurrence of certain predetermined conditions (i.e., NVRAM for persistent
4 storage of file server operations fills up, a preset number of buffers or blocks are cached, etc.).
5

6 In step S706, file server 3 checks flags 7 for each block to be written to storage 4.
7 File server 3 checks flags 7 to see if reservation according to the invention has been activated for
8 the blocks. If reservation has been activated, flow proceeds to step S707. Otherwise, flow skips
9 to step S708.

10 In step S707, reserved block count 13 is decremented by the number of reserved
11 blocks to be written to storage. Because these blocks are actually going to be written to storage
12 4, space for the blocks no longer needs to be reserved. Decrementing the reserved block count
13 releases the reservation of these blocks.
14

15 In step S708, a block or blocks of storage 4 are allocated for the data, and in step
16 S709, the data is written to the block or blocks of storage 4.
17

18 Fig. 8 is a flowchart for explaining a reservation operation for a file for which
19 reservation has already been performed, in which the reservation operation specifies a new file
20 size different from a current (old) file size for the file. For example, the steps in Fig. 8 are
21 performed when a first zero length write specifies a first file size, and then a second zero length
22 write specifies a second different file size.
23

24 In step S801, the current file size is compared with a new file size for the
25 reservation command. If the new file size is less than the old file size, flow proceeds to step
26 S802, where remaining block reservations for the file are released. In order to release the block

1 reservations, flag 15 in inode 14 for the file is cleared, and reserved block count 13 is
2 decremented by the number of blocks still reserved for the file. This number is computed along
3 the lines of step S302 shown in Figs. 3 and 4, based on the current file size, the number of blocks
4 allocated to the file, and the number of cached unallocated blocks for the file.

5
6 If the new file size is greater than the old file size, flow proceeds to step S803,
7 where more blocks are reserved. The number of additional blocks reserved preferably is equal to
8 the difference between the total number of direct and indirect blocks required by the new file size
9 and the total number of direct and indirect blocks required by the current file size. These
10 additional blocks are reserved by incrementing reserved block count 13 by this difference.

11 *Alternative Embodiments*

12
13
14 Although preferred embodiments of the invention are disclosed herein, many
15 variations are possible which remain within the content, scope and spirit of the invention, and
16 these variations would become clear to those skilled in the art after perusal of this application.
17 For example, the invention is equally applicable to file servers and file systems that use different
18 block sizes and block organizations, different buffer cache schemes (e.g., no delayed allocation),
19 different layouts than the WAFL layout, and the like. In addition, the invention is not limited to
20 the particular orderings of steps described above. Therefore, the scope of the invention
21 encompasses the following claims and their legal equivalents and is not limited to the
22 embodiment discussed above.

23 //

24 //

25 //

26 //